

# QuickHeapsort: Modifications and improved analysis

Volker Diekert      Armin Weiß

Universität Stuttgart, FMI  
Universitätsstraße 38  
D-70569 Stuttgart, Germany

`{diekert,weiss}@fmi.uni-stuttgart.de`

September 20, 2012

## Abstract

We present a new analysis for QuickHeapsort splitting it into the analysis of the partition-phases and the analysis of the heap-phases. This enables us to consider samples of non-constant size for the pivot selection and leads to better theoretical bounds for the algorithm.

Furthermore we introduce some modifications of QuickHeapsort, both in-place and using  $n$  extra bits. We show that on every input the expected number of comparisons is  $n \lg n - 0.03n + o(n)$  (in-place) respectively  $n \lg n - 0.997n + o(n)$ . Both estimates improve the previously known best results. (It is conjectured [19] that the in-place algorithm Bottom-Up-Heapsort uses at most  $n \lg n + 0.4n$  on average and for Weak-Heapsort which uses  $n$  extra-bits the average number of comparisons is at most  $n \lg n - 0.42n$  [8].) Moreover, our non-in-place variant can even compete with index based Heapsort variants (e.g. Rank-Heapsort [17]) and Relaxed-Weak-Heapsort ( $n \lg n - 0.9n + o(n)$  comparisons in the worst case) for which no  $\mathcal{O}(n)$ -bound on the number of extra bits is known.

**Keywords.** In-place sorting - heapsort - quicksort - analysis of algorithms

**ACM classification:** F.2.2 Nonnumerical Algorithms and Problems

## 1 Introduction

QuickHeapsort is a combination of Quicksort and Heapsort which was first described by Cantone and Cincotti [2]. It is based on Katajainen's idea for Ultimate Heapsort [12]. In contrast to Ultimate Heapsort it does not have any  $\mathcal{O}(n \lg n)$  bound for the worst case running time. Its advantage is that it is very fast in the average case and hence not only of theoretical interest.

Both algorithms have in common that first the array is partitioned into two parts. Then in one part a heap is constructed and the elements are successively extracted. Finally the remaining elements are treated recursively. The main

advantage of this method is that for the sift-down only one comparison per level is needed, whereas standard Heapsort needs two comparisons per level (for a description of standard Heapsort see some standard textbook, e.g. [6]). This is a severe drawback and one of the reasons why standard Heapsort cannot compete with Quicksort in practice (of course there are also other reasons like cache behavior). Over the time a lot of solutions to this problem appeared like Bottom-Up-Heapsort [19] or MDR-Heapsort [15, 18], which both perform the sift-down by first going down to some leaf and then searching upward for the correct position. Since one can expect that the final position of some introduced element is near to some leaf this is a good heuristic and it leads to provably good results.

The difference between QuickHeapsort and Ultimate Heapsort lies in the choice of the pivot element for partitioning the array. While for Ultimate Heapsort the pivot is chosen as median of the whole array, for QuickHeapsort the pivot is selected as median of some smaller sample (e.g. as median of 3 elements) or by some other method.

In [2] the basic version with fixed index as pivot is analyzed and – together with the median of three version – implemented and compared with other Quick- and Heapsort variants. In [8] Edelkamp and Stiegeler compare these variants with so called Weak-Heapsort [7], some modifications of it (e.g. Relaxed-Weak-Heapsort) and Quick-Weak-Heapsort which relies on the same idea as Quick-Heapsort but uses Weak-Heaps instead of normal heaps. Weak-Heapsort and Quick-Weak-Heapsort beat basic QuickHeapsort with respect to number of comparisons, however they need  $\mathcal{O}(n)$  bits extra-space (for Relaxed-Weak-Heapsort this bound is only conjectured), hence are not in place.

We split the analysis of QuickHeapsort into three parts: the partitioning phases, the heap construction and the heap extraction. This allows us to get better bounds for the running time, especially when choosing the pivot as median of a larger sample. It also simplifies the analysis. We introduce some modifications of QuickHeapsort, too. The first one is in-place and needs  $n \lg n - 0.03n + o(n)$  comparisons on average what is to the best of our knowledge better than any other known in-place Heap- and Quicksort variant. We also examine a modification using  $\mathcal{O}(n)$  bits extra-space, which applies the ideas of MDR-Heapsort to QuickHeapsort. With this method we can bound the average number of comparisons to  $n \lg n - 0.997n + o(n)$ . Actually, a complicated, iterated in-place MergeInsertion uses only  $n \lg n - 1.3n + \mathcal{O}(\lg n)$  comparisons, [16]. Unfortunately, for practical purposes this algorithm is not competitive.

Our contributions are as follows: 1. We give a simplified analysis which gives better bounds than previously known. 2. Our approach yields the first precise analysis of QuickHeapsort when the pivot element is taken from a larger sample. 3. We give a simple in-place modification of QuickHeapsort which saves  $0.75n$  comparisons. 4. We give a modification of QuickHeapsort using  $n$  extra bits only and we can bound the expected number of comparisons. This bound is better than the previously known for the worst case of Heapsort variants using  $\mathcal{O}(n \lg n)$  extra bits for which best and worst case are almost the same. 5. We have implemented QuickHeapsort, and our experiments confirm the theoretical predictions.

The paper is organized as follows: Sect. 2 briefly describes the basic QuickHeapsort algorithm together with our first improvement. In Sect. 3 we analyze the expected running time of QuickHeapsort. Then we introduce some improve-

ments in Sect. 4 allowing  $\mathcal{O}(n)$  additional bits. Finally, in Sect. 5, we present our experimental results comparing the different versions of QuickHeapsort with other Quicksort and Heapsort variants.

## 2 QuickHeapsort

A *two-layer-min-heap* is an array  $A[1..n]$  of  $n$  elements together with a partition  $(G, R)$  of  $\{1, \dots, n\}$  into *green* and *red* elements such that for all  $g \in G, r \in R$  we have  $A[g] \leq A[r]$ . Furthermore, the green elements  $g$  satisfy the heap condition  $A[g] \leq \min\{A[2g], A[2g+1]\}$ , and if  $g$  is red, then  $2g$  and  $2g+1$  are red, too. (The conditions are required to hold, only if the indices involved are in the range of 1 to  $n$ .) The green elements are called “green” because they can be extracted out of the heap without caution, whereas the “red” elements are blocked. *Two-layer-max-heaps* are defined analogously. We can think of a two-layer-heap as rooted binary tree such that each node is either green or red. Green nodes satisfy the standard heap-condition, children of red nodes are red. Two-layer-heaps were defined in [12]. In [2] for the same concept a different language is used (they describe the algorithm in terms of External Heapsort). Now we are ready to describe the QuickHeapsort algorithm as it has been proposed in [2]. Most of it also can be found in pseudocode in App. E.

We intend to sort an array  $A[1..n]$ . First, we choose a *pivot*  $p$ . This is the randomized part of the algorithm. Then, just as in Quicksort, we rearrange the array according to  $p$ . That means, using  $n - 1$  comparisons the partitioning function returns an index  $k$  and rearranges the array  $A$  so that  $A[i] \geq A[k]$  for  $i < k$ ,  $A[k] = p$ , and  $A[k] \geq A[j]$  for  $k < j$ . After the partitioning a two-layer-heap is built out of the elements of the smaller part of the array, either the part left of the pivot or right of the pivot. We call this smaller part *heap-area* and the larger part *work-area*. More precisely, if  $k - 1 < n - k$ , then  $\{1, \dots, k - 1\}$  is the heap-area and  $\{k + 1, \dots, n\}$  is the work-area. If  $k - 1 \geq n - k$ , then  $\{1, \dots, k - 1\}$  is the work-area and  $\{k + 1, \dots, n\}$  is the heap-area. Note that we know the final position of the pivot element without any further comparison. Therefore, we do not count it to the heap-area nor to the work-area. If the heap-area the part of the array left of the pivot, a two-layer-max-heap is built, otherwise a two-layer-min-heap is built.

At the beginning the heap-area is an ordinary heap, hence it is a two-layer-heap consisting of green elements, only. Now the heap extraction phase starts. Let  $m$  denote the size of the heap-area. If we are in the case of a max-heap, these  $m$  elements are moved to the back of the array, in the case of a min-heap they are moved to the front of the array (which is in both cases the work-area). For a max-heap, the extraction of one element works as follows: the root of the heap is placed at the current position of the work-area (which at the beginning is its last position). Then, starting from the root the resulting “hole” is trickled down: always the larger child is moved up into the vacant position and then this child is treated recursively. This stops as soon as a leaf is reached. We call this the *SpecialLeaf* procedure (Alg. 5.2) according to [2]. Now, the element which before was at the current position in the work-area is placed as red element in this hole at the leaf in the heap-area. Finally the current position in the work-area is moved by one and the next element can be extracted.

The procedure sorts correctly, because after the partitioning it is guaranteed

that all red elements are smaller than all green elements. Furthermore there is enough space in the work-area to place all green elements of the heap, since the heap is always the smaller part of the array. After extracting all green elements the pivot element is placed at its final position and the remaining elements are sorted recursively.

Actually we can improve the procedure, thereby saving  $3n/4$  comparisons by a simple trick. Before the heap extraction phase starts in the heap-area with  $m$  elements, we perform at most  $\frac{m+2}{4}$  additional comparisons in order to arrange all pairs of leaves which share a parent such that the left child is not smaller than its right sibling. Now, in every call of `SpecialLeaf`, we can save exactly one comparison, since we do not need to compare two leaves. For a max-heap we only need to move up the left child and put the right one at the place of the former left one. Summing up over all heaps during an execution of standard `QuickHeapsort`, we invest  $\frac{n+2t}{4}$  comparisons in order to save  $n$  comparisons, where  $t$  is the number of recursive calls. The expected number of  $t$  is in  $\mathcal{O}(\lg n)$ . Hence, we can expect to save  $\frac{3n}{4} + \mathcal{O}(\lg n)$  comparisons. We call this version the *improved* variant of `QuickHeapsort`.

### 3 Analysis of QuickHeapsort

This section contains the main contribution of the paper. We analyze the number of comparisons.

Throughout the logarithm  $\lg$  is always meant to base 2. By  $n$  we denote the number of elements of an array to be sorted. We use standard  $\mathcal{O}$ -notation where  $\mathcal{O}(g)$ ,  $o(g)$ , and  $\omega(g)$  denote classes of functions. In our analysis we do not assume any random distribution of the input, i.e. it is valid for every permutation of the input array. Randomization is used however for pivot selection. With  $\Pr[e]$  we denote the probability of some event  $e$ . The expected value of a random variable  $T$  is denoted by  $\mathbb{E}[T]$ .

The number of assignments is bounded by some small constant times the number of comparisons. Let  $T(n)$  denote the number of comparisons during `QuickHeapsort` on a fixed array of  $n$  elements. We are going to split the analysis of `QuickHeapsort` into three parts:

1. Partitioning with an expected number of comparisons  $\mathbb{E}[T_{\text{part}}(n)]$  (average case).
2. Heap construction with at most  $T_{\text{con}}(n)$  comparisons (worst case).
3. Heap extraction (sorting phase) with at most  $T_{\text{ext}}(n)$  comparisons (worst case).

We analyze the three parts separately and put them together at the end. The partitioning is the only randomized part of our algorithm. The expected number of comparisons depends on the selection method for the pivot. For the expected number of comparisons by `QuickHeapsort` on the input array we obtain  $\mathbb{E}[T(n)] \leq T_{\text{con}}(n) + T_{\text{ext}}(n) + \mathbb{E}[T_{\text{part}}(n)]$ .

**Theorem 3.1** *Let  $f \in \omega(1) \cap o(n)$  with  $1 \leq f(n) \leq n$ , e.g.,  $f(n) = \lg n$ , and let  $\mathbb{E}[T(n)]$  be the expected number of comparisons by `QuickHeapsort` on a fixed*

input array of size  $n$ . Choosing the pivot as median of  $f(n)$  randomly selected elements in time  $\mathcal{O}(f(n))$ , we have

$$\begin{aligned}\mathbb{E}[T(n)] &\leq n \lg n + 0.72n + o(n) \\ \mathbb{E}[T(n)] &\leq n \lg n - 0.03n + o(n) \text{ improved variant.}\end{aligned}$$

The significance of the selection method for the pivot becomes evident, if we compare Thm. 3.1 with bounds one can derive for simpler selection methods as used in Prop. 3.2.

**Proposition 3.2** *Choosing the pivot a.) at random, respectively b.) as median-of-three, we have the following bounds*

$$\begin{aligned}a.) \mathbb{E}[T(n)] &\leq n \lg n + 2.72n + o(n), \quad \mathbb{E}[T(n)] \leq n \lg n + 1.97n + o(n) \text{ impr.} \\ b.) \mathbb{E}[T(n)] &\leq n \lg n + 1.92n + o(n), \quad \mathbb{E}[T(n)] \leq n \lg n + 1.17n + o(n) \text{ impr.}\end{aligned}$$

The proof of these results are postponed to Sect. 3.3. Note that it is enough to prove the results without the improvement, since the difference is always  $0.75n$ .

### 3.1 Heap Construction

The standard heap construction [9] needs at most  $2m$  comparisons to construct a heap of size  $m$  in the worst case and approximately  $1.88m$  in the average case. For the mathematical analysis better theoretical bounds can be used. The best result we are aware of is due to Chen et al. in [5]. According to this result we have  $T_{\text{con}}(m) \leq 1.625m + o(m)$ .

Earlier results are of similar magnitude, by [4] it has been known that  $T_{\text{con}}(m) \leq 1.632m + o(m)$  and by [10] it has been known  $T_{\text{con}}(m) \leq 1.625m + o(m)$ , but Gonnet and Munro used  $\mathcal{O}(m)$  extra bits to get this result, whereas the new result of Chen et al. is in-place (by using only  $\mathcal{O}(\lg m)$  extra bits).

During the execution of QuickHeapsort over  $n$  elements, every element is part of a heap only once. Hence, the sizes of all heaps during the entire procedure sum up to  $n$ . With the result of [5] the total number of comparisons performed in the construction of all heaps satisfies:

**Proposition 3.3**

$$T_{\text{con}}(n) \leq 1.625n + o(n).$$

### 3.2 Heap Extraction

For a real number  $r \in \mathbb{R}$  with  $r > 0$  we define  $\{r\}$  by the following condition

$$r = 2^k + \{r\} \text{ with } k \in \mathbb{Z} \text{ and } 0 \leq \{r\} < 2^k.$$

This means that  $2^k$  is largest power of 2 which is less than or equal to  $r$  and  $\{r\}$  is the difference to that power, i.e.  $\{r\} = r - 2^{\lfloor \lg r \rfloor}$ . In this section we first analyze the extraction phase of one two-layer-heap of size  $m$ . After that, we bound the number of comparisons  $T_{\text{ext}}(n)$  performed in the worst case during all heap extraction phases of one execution of QuickHeapsort on an array of size  $n$ . Thm. 3.4 is our central result about heap extraction.

### Theorem 3.4

$$T_{\text{ext}}(n) \leq n \cdot (\lfloor \lg n \rfloor - 3) + 2\{n\} + \mathcal{O}(\lg^2 n).$$

The proof of Thm. 3.4 covers the almost rest of Section 3.2. In the following, the *height*  $\text{height}(v)$  of an element  $v$  in a heap  $H$  is the maximal distance from the that node to a leaf below it. The *height* of  $H$  is the height of its root. The *level*  $\text{level}(v)$  of  $v$  to be its distance from the root.

In this section we want to count the comparisons during SpecialLeaf procedures, only. Recall that a SpecialLeaf procedure is a cyclic shift on a path from the root down to some leaf, and the number comparisons is exactly the length of this path. Hence the upper bound is the height of the heap. But there is a better analysis.

Let us consider a heap with  $m$  green elements which are all extracted by SpecialLeaf procedures. The picture is as follows: First, we color the green root red. Next, we perform a cyclic shift defined by the SpecialLeaf procedure. In particular, the leaf is now red. Moreover, red positions remain red, but there is exactly one position  $v$  which has changed its color from green to red. This position  $v$  is on the path defined by the SpecialLeaf procedure. Hence, the number of comparisons needed to color the position  $v$  red is bounded by  $\text{height}(v) + \text{level}(v)$ .

The total number of comparisons  $E(m)$  to extract all  $m$  elements of a Heap  $H$  is therefore bounded by

$$E(m) \leq \sum_{v \in H} (\text{height}(v) + \text{level}(v)). \quad (1)$$

We have  $\text{height}(H) - 1 \leq \text{height}(v) + \text{level}(v) \leq \text{height}(H) = \lfloor \lg m \rfloor$  for all  $v \in H$ . We now count the number of elements  $v$  where  $\text{height}(v) + \text{level}(v) = \lfloor \lg m \rfloor$  and the number of elements  $v$  where  $\text{height}(v) + \text{level}(v) = \lfloor \lg m \rfloor - 1$ . Since there are exactly  $\{m\} + 1$  nodes of level  $\lfloor \lg m \rfloor$ , there are at most  $2\{m\} + 1 + \lg m$  elements  $v$  with  $\text{height}(v) + \text{level}(v) = \lfloor \lg m \rfloor$ . All other elements satisfy  $\text{height}(v) + \text{level}(v) = \lfloor \lg m \rfloor - 1$ . We obtain

$$\begin{aligned} E(m) &\leq 2 \cdot \{m\} \cdot \lfloor \lg m \rfloor + (m - 2 \cdot \{m\})(\lfloor \lg m \rfloor - 1) + \mathcal{O}(\lg m) \\ &= m \cdot (\lfloor \lg m \rfloor - 1) + 2 \cdot \{m\} + \mathcal{O}(\lg m). \end{aligned} \quad (2)$$

Note that this is an estimate of the worst case, however this analysis also shows that the best case only differs by  $\mathcal{O}(\lg m)$ -terms from the worst case.

Now, we want to estimate the number of comparisons in the worst case performed during all heap extraction phases together. During QuickHeapsort over  $n$  elements we create a sequence  $H_1, \dots, H_t$  of heaps of green elements which are extracted using the SpecialLeaf procedure. Let  $m_i = |H_i|$  be the size of the  $i$ -th Heap. The sequence satisfies  $2m_i \leq n - \sum_{j < i} m_j$ , because heaps are constructed and extracted on the smaller part of the array.

Here comes a subtle observation: Assume that  $m_1 + m_2 \leq n/2$ . If we replace the first two heaps with one heap  $H'$  of size  $|H'| = m_1 + m_2$ , then the analysis using the sequence  $H', H_3, \dots, H_t$  cannot lead to a better bound. Continuing this way, we may assume that we have  $t \in \mathcal{O}(\lg n)$  and therefore

$\sum_{1 \leq i \leq t} \mathcal{O}(\lg m_i) \subseteq \mathcal{O}(\lg^2 n)$ . With Eq. (2) we obtain the bound

$$T_{\text{ext}}(n) \leq \sum_{i=1}^t E(m_i) = \left( \sum_{i=1}^t m_i \cdot \lfloor \lg m_i \rfloor + 2 \{m_i\} \right) - n + \mathcal{O}(\lg^2 n). \quad (3)$$

Later we will replace the  $m_i$  by other positive real numbers. Therefore we define the following notion. Let  $1 \leq \nu \in \mathbb{R}$ . We say a sequence  $x_1, x_2, \dots, x_t$  with  $x_i \in \mathbb{R}^{>0}$  is *valid* w.r.t.  $\nu$ , if for all  $1 \leq i \leq t$  we have:

$$2x_i \leq \nu - \sum_{j < i} x_j.$$

As just mentioned the initial sequence  $m_1, m_2, \dots, m_t$  is valid w.r.t.  $n$ . Let us define a continuous function  $F : \mathbb{R}^{>0} \rightarrow \mathbb{R}$  by

$$F(x) = x \cdot \lfloor \lg x \rfloor + 2 \{x\}.$$

It is continuous, since for  $x = 2^k$ ,  $k \in \mathbb{Z}$  we have  $F(x) = xk = \lim_{\varepsilon \rightarrow 0} (x - \varepsilon)(k-1) + 2 \{x - \varepsilon\}$ . It is piecewise differentiable with right derivative  $\lfloor \lg x \rfloor + 2$ . Therefore:

**Lemma 3.5** *Let  $x \geq y > \delta \geq 0$ . Then we have the inequalities:*

$$F(x) + F(y) \leq F(x + \delta) + F(y - \delta) \text{ and } F(x) + F(y) \leq F(x + y).$$

**Lemma 3.6** *Let  $1 \leq \nu \in \mathbb{R}$ . For all sequences  $x_1, x_2, \dots, x_t$  with  $x_i \in \mathbb{R}^{>0}$ , which are valid w.r.t.  $\nu$ , we have:*

$$\sum_{i=1}^t F(x_i) \leq \sum_{i=1}^{\lfloor \lg \nu \rfloor} F\left(\frac{\nu}{2^i}\right).$$

**Proof.** The result is true for  $\nu \leq 2$ , because then  $F(x_i) \leq F(\nu/2) \leq F(1) = 0$  for all  $i$ . Thus, we may assume  $\nu \geq 2$ . We perform induction on  $t$ . For  $t = 1$  the statement is clear, since  $\lg \nu \geq 1$  and  $x_1 \leq \nu/2$ . Now let  $t > 1$ . By Lem. 3.5, we have  $F(x_1) + F(x_2) < F(x_1 + x_2)$ . Now, if  $x_1 + x_2 \leq \frac{\nu}{2}$ , then the sequence  $x_1 + x_2, x_3, \dots, x_t$  is valid, too; and we are done by induction. Hence, we may assume  $x_1 + x_2 > \frac{\nu}{2}$ . If  $x_1 \leq x_2$ , then

$$2x_1 = 2x_2 + 2(x_1 - x_2) \leq \nu - x_1 + 2(x_1 - x_2) = \nu - x_2 + x_1 - x_2 \leq \nu - x_2.$$

Thus, if  $x_1 \leq x_2$ , then the sequence  $x_2, x_1, x_3, \dots, x_t$  is valid, too. Thus, it is enough to consider  $x_1 \geq x_2$  with  $x_1 + x_2 > \frac{\nu}{2}$ .

We have  $\frac{\nu}{2} \geq 1$  and the sequence  $x'_2, x_3, \dots, x_t$  with  $x'_2 = x_1 + x_2 - \frac{\nu}{2}$  is valid w.r.t.  $\nu/2$ , because

$$x'_2 = x_1 + x_2 - \frac{\nu}{2} \leq x_1 + \frac{\nu - x_1}{2} - \frac{\nu}{2} = \frac{x_1}{2} \leq \frac{\nu}{4}.$$

Therefore, by induction on  $t$  and Lem. 3.5 we obtain the claim:

$$\sum_{i=1}^t F(x_i) \leq F(\nu/2) + F(x'_2) + \sum_{i=3}^t F(x_i) \leq F(\nu/2) + \sum_{i=2}^{\lfloor \lg \nu \rfloor} F\left(\frac{\nu}{2^i}\right) \leq \sum_{i=1}^{\lfloor \lg \nu \rfloor} F\left(\frac{\nu}{2^i}\right).$$

□

**Lemma 3.7**

$$\sum_{i=1}^{\lfloor \lg n \rfloor} F\left(\frac{n}{2^i}\right) \leq F(n) - 2n + \mathcal{O}(\lg n).$$

**Proof.**

$$\begin{aligned} \sum_{i=1}^{\lfloor \lg n \rfloor} F\left(\frac{n}{2^i}\right) &= n \lfloor \lg n \rfloor \cdot \sum_{i=1}^{\lfloor \lg n \rfloor} \frac{1}{2^i} - n \cdot \sum_{i=1}^{\lfloor \lg n \rfloor} \frac{i}{2^i} + 2 \{n\} \cdot \sum_{i=1}^{\lfloor \lg n \rfloor} \frac{1}{2^i} \\ &\leq n \lfloor \lg n \rfloor \cdot \sum_{i \geq 1} \frac{1}{2^i} - n \cdot \sum_{i \geq 1} \frac{i}{2^i} + 2 \{n\} \cdot \sum_{i \geq 1} \frac{1}{2^i} + \frac{n}{2^{\lfloor \lg n \rfloor}} \cdot \sum_{i > 0} \frac{i + \lfloor \lg n \rfloor}{2^i} \\ &= n \lfloor \lg n \rfloor - 2n + 2 \{n\} + \mathcal{O}(\lg n). \end{aligned}$$

□

Applying these lemmata to Eq. (3) yields the proof of Thm. 3.4.

**Corollary 3.8** *We have*

$$T_{\text{ext}}(n) \leq n \lg n - 2.9139n + \mathcal{O}(\lg^2 n)$$

**Proof.** By [18, Thm. 1] we have  $F(n) - 2n \leq n \lg n - 1.9139n$ . Hence, Cor. 3.8 follows directly from Thm. 3.4. □

### 3.3 Partitioning

In the following  $T_{\text{pivot}}(n)$  denotes the number of comparisons required to choose the pivot element in the worst case; and, as before,  $\mathbb{E}[T_{\text{part}}(n)]$  denotes the expected number of comparisons performed during partitioning. We have the following recurrence:

$$\mathbb{E}[T_{\text{part}}(n)] \leq n - 1 + T_{\text{pivot}}(n) + \sum_{k=1}^n \Pr[\text{pivot} = k] \cdot \mathbb{E}[T_{\text{part}}(\max\{k-1, n-k\})]. \quad (4)$$

If we choose the pivot at random, then we obtain by standard methods:

$$\mathbb{E}[T_{\text{part}}(n)] \leq n - 1 + \frac{1}{n} \cdot \sum_{k=1}^n \mathbb{E}[T_{\text{part}}(\max\{k-1, n-k\})] \leq 4n. \quad (5)$$

Similarly, if we choose the pivot with the median-of-three, then we obtain:

$$\mathbb{E}[T_{\text{part}}(n)] \leq 3.2n + \mathcal{O}(\lg n). \quad (6)$$

The proof of Prop. 3.2 follows from Equations (5) and (6), Thm. 3.4, and Prop. 3.3. Using a growing number of elements (as  $n$  grows) as sample for the pivot selection, we can do better. Thm. 3.1 follows now from Thm. 3.4, Prop. 3.3, and Thm. 3.9.



**Theorem 3.9** *Let  $f \in \omega(1) \cap o(n)$  with  $1 \leq f(n) \leq n$ . When choosing the pivot as median of  $f(n)$  randomly selected elements in time  $\mathcal{O}(f(n))$  (e.g. with the algorithm of [1]), the expected number of comparisons used in all recursive calls of partitioning is in  $2n + o(n)$ .*

Thm. 3.9 is close to a well-known result in [14, Thm. 5] on Quickselect, see Cor. 3.11. Formally speaking we cannot use it directly, because we deal with QuickHeapsort, where after partitioning the recursive call is on the larger part. Because of that, and for the sake of completeness, we give a proof. Moreover, our proof is elementary and simpler than the one in [14]. The key point is Lem. 3.10. Its proof is rather standard, see App. B for details.

**Lemma 3.10** *Let  $0 < \delta < \frac{1}{2}$ . If we choose the pivot as median of  $2c + 1$  elements such that  $2c + 1 \leq \frac{n}{2}$ , then we have  $\Pr[\text{pivot} \leq \frac{n}{2} - \delta n] < (2c + 1)\alpha^c$  where  $\alpha = 4(\frac{1}{4} - \delta^2) < 1$ .*

**Proof of Thm. 3.9.** As an abbreviation, we let  $E(n) = \mathbb{E}[T_{\text{part}}(n)]$  be the expected number of comparisons performed during partitioning. We are going to show that for all  $\epsilon > 0$  there is some  $D \in \mathbb{R}$  such that

$$E(n) < (2 + \epsilon)n + D. \quad (7)$$

So, we fix some  $1 \geq \epsilon > 0$ . We choose  $\delta > 0$  such that  $(2 + \epsilon)\delta < \frac{\epsilon}{4}$ . Moreover, for this proof let  $\mu = \frac{n+1}{2}$ . Positions of possible pivots  $k$  with  $\mu - \delta n \leq k \leq \mu + \delta n$  form a small fraction of all positions, and they are located around the median. Nevertheless, applying Lem. 3.10 with  $c = f(n) \in \omega(1) \cap o(n)$  yields for all  $n$ , which are large enough:

$$\Pr[\text{pivot} < \mu - \delta n] \leq (2f(n) + 1) \cdot \alpha^{f(n)} \leq \frac{1}{48}\epsilon. \quad (8)$$

The analogous inequality holds for  $\Pr[\text{pivot} > \mu + \delta n]$ . Because  $T_{\text{pivot}}(n) \in o(n)$ , we have

$$T_{\text{pivot}}(n) \leq \frac{1}{8}\epsilon n. \quad (9)$$

for  $n$  large enough. Now, we choose  $n_0$  such that Eq. (8) and Eq. (9) hold for  $n \geq n_0$  and such that we have  $(2 + \epsilon)\delta + \frac{2}{n_0} < \frac{\epsilon}{4}$ . We set  $D = E(n_0) + 1$ . Hence for  $n < n_0$  the desired result Eq. (7) holds. Now, let  $n \geq n_0$ . From Eq. (4) we obtain by symmetry:

$$\begin{aligned} E(n) &\leq n - 1 + T_{\text{pivot}}(n) \\ &\quad + \sum_{k=\lceil \mu - \delta n \rceil}^{\lfloor \mu + \delta n \rfloor} \Pr[\text{pivot} = k] \cdot E(k - 1) \\ &\quad + 2 \sum_{k=\lfloor \mu + \delta n \rfloor + 1}^n \Pr[\text{pivot} = k] \cdot E(k - 1). \end{aligned}$$

Since  $E$  is monotone,  $E(k)$  can be bounded by the highest value in the respective interval:

$$\begin{aligned}
&\leq n + \frac{1}{8}\epsilon n \\
&\quad + \Pr[\mu - \delta n \leq \text{pivot} \leq \mu + \delta n] \cdot E(\lfloor \mu + \delta n \rfloor) \\
&\quad + 2 \Pr[\text{pivot} > \mu + \delta n] \cdot E(n-1) \\
&\leq n + \frac{1}{8}\epsilon n + \left(1 - \frac{1}{24}\epsilon\right) \cdot E(\lfloor \mu + \delta n \rfloor) + 2 \frac{1}{48}\epsilon \cdot E(n-1).
\end{aligned}$$

By induction we assume  $E(k) \leq (2 + \epsilon)k + D$  for  $k < n$ . Hence:

$$\begin{aligned}
E(n) &\leq n + \frac{1}{8}\epsilon n + \left(1 - \frac{1}{24}\epsilon\right) \cdot ((2 + \epsilon) \cdot (\mu + \delta n) + D) + \frac{1}{24}\epsilon \cdot ((2 + \epsilon)n + D) \\
&\leq n + (2 + \epsilon) \cdot \left(\frac{n+1}{2} + \delta n\right) + \frac{1}{8}\epsilon n + \frac{1}{24}\epsilon(2 + \epsilon)n + D \\
&\leq 2n + 1 + \frac{\epsilon}{2} + (2 + \epsilon)\delta n + \frac{3}{4}\epsilon n + D < (2 + \epsilon)n + D.
\end{aligned}$$

□

**Corollary 3.11 ([14])** *Let  $f \in \omega(1) \cap o(n)$  with  $1 \leq f(n) \leq n$ . When implementing Quickselect with the median of  $f(n)$  randomly selected elements as pivot, the expected number of comparisons is  $2n + o(n)$ .*

**Proof.** In QuickHeapsort the recursion is always on the larger part of the array. Hence, the number of comparisons in partitioning for QuickHeapsort is an upper bound on the number of comparisons in Quickselect. □

In [14] it is also proved that choosing the pivot as median of  $\mathcal{O}(\sqrt{n})$  elements is optimal for Quicksort as well as for Quickselect. This suggests that we choose the same value in QuickHeapsort; what is backed by our experiments.

## 4 Modifications of QuickHeapsort using Extra-space

In this section we want to describe some modification of QuickHeapsort using  $n$  bits of extra storage. We introduce two bit-arrays. In one of them (the CompareArray) – which is actually two bits per element – we store the comparisons already done (we need two bits, because there are three possible values – right, left, unknown – we have to store). In the other one (the RedGreenArray) we store which element is red and which is green.

Since the heaps have maximum size  $n/2$ , the RedGreenArray only requires  $n/2$  bits. The CompareArray is only needed for the inner nodes of the heaps, i.e. length  $n/4$  is sufficient. Totally this sums up to  $n$  extra bits.

For the heap construction we do not use the algorithms described in Sect. 3.1. With the CompareArray we can do better by using the algorithm of McDiarmid and Reed [15].

The heap construction works similarly to Bottom-Up-Heapsort, i.e. the array is traversed backward calling for all inner positions  $i$  the Reheap procedure on  $i$ . The Reheap procedure takes the subheap with root  $i$  and restores the heap condition, if it is violated at the position  $i$ . First, the Reheap procedure determines a *special leaf* using the SpecialLeaf procedure as described in Sect. 2, but without moving the elements. Then, the final position of the former root is determined going upward from the special leaf (bottom-up-phase). At the end, the elements above this final position are moved up towards the root by one position. That means that all but one element which are compared during the bottom-up-phase, stay in their places. Since in the SpecialLeaf procedure these elements have been compared with their siblings, these comparisons can be stored in the CompareArray and can be used later.

With another improvement concerning the construction of heaps with seven elements as in [3] the benefits of this array can be exploited even more.

The RedGreenArray is used during the sorting phase, only. Its functionality is straightforward: Every time a red element is inserted into the heap, the corresponding bit is set to red. The SpecialLeaf procedure can stop as soon as it reaches an element without green children. Whenever a red and a green element have to be compared, the comparison can be skipped.

**Theorem 4.1** *Let  $f \in \omega(1) \cap o(n)$  with  $1 \leq f(n) \leq n$ , e.g.,  $f(n) = \lg n$ , and let  $\mathbb{E}[T(n)]$  be the expected number of comparisons by QuickHeapsort using the CompareArray with the improvement of [3] and the RedGreenArray on a fixed input array of size  $n$ . Choosing the pivot as median of  $f(n)$  randomly selected elements in time  $\mathcal{O}(f(n))$ , we have*

$$\mathbb{E}[T(n)] \leq n \lg n - 0.997n + o(n). \quad (1)$$

We can analyze the savings by the two arrays separately, because the CompareArray only affects comparisons between two green elements, while the RedGreenArray only affects comparisons involving at least one red element. First, we consider the heap construction using the CompareArray. With this array we obtain the same worst case bound as for the standard heap construction method. However, the CompareArray has the advantage that at the end of the heap construction many comparisons are stored in the array and can be reused for the extraction phase. More precisely: For every comparison except the first one made when going upward from the special leaf, one comparison is stored in the CompareArray, since for every additional comparison one element on the path defined by SpecialLeaf stays at its place. Because every pair of siblings has to be compared at one point during the heap construction or extraction, all these stored comparisons can be reused. Hence, we only have to count the comparisons in the SpecialLeaf procedure during the construction plus  $\frac{n}{2}$  for the first comparison when going upward. Thus, we get an amortized bound for the comparisons during construction of  $\frac{3n}{2}$ .

In [3] the notion of *Fine-Heaps* is introduced. A Fine Heap is a heap with the additional CompareArray such that for every node the larger child is stored in the array. Such a Fine-Heap of size  $m$  can be constructed using the above method with  $2m$  comparisons. In [3] Carlsson, Chen and Mattsson showed that a Fine-Heap of size  $m$  actually can be constructed with only  $\frac{23}{12}m + \mathcal{O}(\lg^2 m)$  comparisons. That means we have to invest  $\frac{23}{12}m + \mathcal{O}(\lg^2 m)$  for the heap construction and at the end there are  $\frac{m}{2}$  comparisons stored in the array. All these

comparisons stored in the array are used later. Summing up over all heaps during an execution of QuickHeapsort, we can save another  $\frac{1}{12}n$  comparisons additionally to the comparisons saved by the CompareArray with the result of [3].

Hence, for the amortized cost of the heap construction  $T_{\text{con}}^{\text{amort}}$  (i.e. the number of comparisons needed to build the heap minus the number of comparisons stored in the CompareArray after the construction which all can be reused later) we have obtained:

**Proposition 4.2**

$$T_{\text{con}}^{\text{amort}}(n) \leq \frac{17}{12}n + o(n).$$

This bound is slightly better than the average case for the heap construction with the algorithm of [15] which is  $1.52n$ .

Now, we want to count the number of comparisons we save using the RedGreenArray. We distinguish the two cases that two red elements are compared and that a red and a green element are compared. Every position in the heap has to turn red at one point. At that time, all nodes below this position are already red. Hence, for that element we save as many comparisons as the element is above the bottom level. Summing over all levels of a heap of size  $m$  the saving results in:

$$\approx \frac{m}{4} \cdot 1 + \frac{m}{8} \cdot 2 + \dots = m \cdot \sum_{i \geq 1} i 2^{-i-1} = m.$$

This estimate is exact up to  $\mathcal{O}(\lg m)$ -terms. Since the expected number of heaps is  $\mathcal{O}(\lg n)$ , we obtain for the overall saving the value  $T_{\text{saveRR}}(n) = n + \mathcal{O}(\lg^2 n)$ .

Another place where we save comparisons with the RedGreenArray is when a red element is compared with a green element. It occurs at least one time – when the node loses its last green child – for every inner node that we compare a red child with a green child. Hence, we save at least as many comparisons as there are inner nodes with two children, i.e. at least  $\frac{m}{2} - 1$ . Since every element – except the expected  $\mathcal{O}(\lg n)$  pivot elements – is part of a heap exactly once, we save at least

$$T_{\text{saveRG}}(m) \geq \frac{m}{2} + \mathcal{O}(\lg n).$$

comparisons when comparing green with red elements. In the average case the saving might be even slightly higher, since comparisons can also be saved when a node does not lose its last green child.

Summing up all our savings and using the median of  $f(n) \in \omega(1) \cap o(n)$  as pivot we obtain the proof of Thm. 4.1:

$$\begin{aligned} \mathbb{E}[T(n)] &\leq T_{\text{con}}^{\text{amort}}(n) + T_{\text{ext}}(n) + \mathbb{E}[T_{\text{part}}(n)] - T_{\text{saveRR}}(n) - T_{\text{saveRG}}(n) \\ &\leq \frac{17}{12}n + n \cdot (\lfloor \lg n \rfloor - 3) + 2\{n\} + 2n - \frac{3n}{2} + o(n) \\ &\leq n \lg n - 0.997n + o(n). \end{aligned}$$

## 5 Experimental Results and Conclusion

In Table 1 we compare the different versions of QuickHeapsort we considered in this paper, i.e. the basic version, the improved variant of Sect. 2, and the version using bit-arrays (however, without the modification by [3]). We do not present actual running times, because the main focus of this paper is on the number of comparisons and therefore our implementation is not tuned to minimize running times (in [11] it is shown that in order to minimize the running time of Quicksort it can be even better to invest some more comparisons by choosing a ‘bad’ pivot – considering such effects would go far beyond the scope of this work). More results with other pivot selection strategies are in Table 2 in App. C confirming that a sample size of  $\sqrt{n}$  is optimal for pivot selection with respect to the number of comparisons and also that the  $o(n)$ -terms in Thm. 3.1 and Thm. 3.9 are not too big. For the heap construction we implemented the normal algorithm due to Floyd [9] as well as the algorithm using the extra bit-array (which is the same as in MDR-Heapsort). All algorithms are implemented with median of 3 and with median of  $\sqrt{n}$  elements as pivot. We compare them with Quicksort implemented with the same pivot selection strategies, Bottom-Up-Heapsort and MDR-Heapsort. In Table 1 we also added the values for Relaxed-Weak-Heapsort which were presented in [8]. All the numbers displayed here are average values over 100 runs with random data. As our theoretical estimates predict, QuickHeapsort with bit-arrays beats all other variants including Relaxed-Weak-Heapsort when implemented with median of  $\sqrt{n}$  for pivot selection (it needs  $260584 \approx 0.26 \cdot 10^6$  comparisons less than Relaxed-Weak-Heapsort). It also performs  $326728 \approx 0.33 \cdot 10^6$  comparisons less than

Sorting algorithm	Average number of comparisons for $n = 10^6$
Basic QuickHeapsort with median of 3	21327478
Basic QuickHeapsort with median of $\sqrt{n}$	20783631
Improved QuickHeapsort with median of 3	20639046
Improved QuickHeapsort with median of $\sqrt{n}$	20135688
QuickHeapsort with bit-arrays with median of 3	19207289
QuickHeapsort with bit-arrays with median of $\sqrt{n}$	18690841 *Best result*
Quicksort with median of 3	21491310
Quicksort with median of $\sqrt{n}$	19548149
Bottom-Up-Heapsort	20294866
MDR-Heapsort	20001084
Relaxed-Weak-Heapsort	18951425
<b>Lower Bound: <math>\lg n!</math></b>	$18488884 \approx \lg(10^6!)$

Table 1: QuickHeapsort and other algorithms tested on  $10^6$  elements (the data for Relaxed-Weak-Heapsort is taken from [8]).

our theoretical predictions which are  $10^6 \cdot \lg(10^6) - 0.9139 \cdot 10^6 \approx 19017569$  comparisons.

In this paper we have shown that with known techniques QuickHeapsort can be implemented with expected number of comparisons less than  $n \lg n -$

$0.03n + o(n)$  and extra storage  $O(1)$ . On the other hand, using  $n$  extra bits we can improve this to  $n \lg n - 0.997n + o(n)$ , i.e. we showed that QuickHeapsort can compete with the most advanced Heapsort variants. These theoretical estimates were also confirmed by our experiments. We also considered different pivot selection schemes. For any constant size sample for pivot selection, QuickHeapsort beats Quicksort for large  $n$ , since Quicksort has a expected running time of  $\approx Cn \lg n$  with  $C > 1$ . However, when choosing the pivot as median of  $\sqrt{n}$  elements (i.e. with the optimal strategy) then our experiments show that Quicksort needs less comparisons than QuickHeapsort. However, using bit-arrays QuickHeapsort is the winner, again. In order to make the last statement rigorous, better theoretical bounds for Quicksort with sampling  $\sqrt{n}$  elements are needed. For future work it would also be of interest to prove the optimality of  $\sqrt{n}$  elements for pivot selection in QuickHeapsort, to estimate the lower order terms of the average running time of QuickHeapsort and also to find an exact average case analysis for the saving by the bit-arrays.

## References

- [1] M. Blum, R. W. Floyd, V. Pratt, R. L. Rivest, and R. E. Tarjan. Time bounds for selection. *J. Comput. Syst. Sci.*, 7(4):448–461, 1973.
- [2] D. Cantone and G. Cincotti. QuickHeapsort, an efficient mix of classical sorting algorithms. *Theor. Comput. Sci.*, 285(1):25–42, 2002.
- [3] S. Carlsson, J. Chen, and C. Mattsson. Heaps with Bits. In D.-Z. Du and X.-S. Zhang, editors, *ISAAC*, volume 834 of *LNCS*, pages 288–296. Springer, 1994.
- [4] J. Chen. A Framework for Constructing Heap-like structures in-place. In K.-W. Ng et al., editors, *ISAAC*, volume 762 of *LNCS*, pages 118–127. Springer, 1993.
- [5] J. Chen, S. Edelkamp, A. Elmasry, and J. Katajainen. In-place Heap Construction with Optimized Comparisons, Moves, and Cache Misses. In B. Rován, V. Sassone, and P. Widmayer, editors, *MFCS*, volume 7464 of *LNCS*, pages 259–270. Springer, 2012.
- [6] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*. The MIT Press, 3 edition, 2009.
- [7] R. D. Dutton. Weak-heap sort. *BIT*, 33(3):372–381, 1993.
- [8] S. Edelkamp and P. Stiegeler. Implementing HEAPSORT with  $n \lg n - 0.9n$  and QUICKSORT with  $n \lg n + 0.2n$  comparisons. *ACM J. of Experimental Algorithmics*, 7:5, 2002.
- [9] R. W. Floyd. Algorithm 245: Treesort. *Commun. ACM*, 7(12):701, 1964.
- [10] G. H. Gonnet and J. I. Munro. Heaps on Heaps. *SIAM J. Comput.*, 15(4):964–971, 1986.

- [11] K. Kaligosi and P. Sanders. How Branch Mispredictions Affect Quicksort. In Y. Azar and T. Erlebach, editors, *ESA*, volume 4168 of *LNCS*, pages 780–791. Springer, 2006.
- [12] J. Katajainen. The Ultimate Heapsort. In X. Lin, editor, *CATS*, volume 20 of *Australian Computer Science Communications*, pages 87–96. Springer-Verlag 1998.
- [13] D. E. Knuth. *The art of computer programming. Vol. 3*. Addison-Wesley, 1998.
- [14] C. Martínez and S. Roura. Optimal Sampling Strategies in Quicksort and Quickselect. *SIAM J. Comput.*, 31(3):683–705, 2001.
- [15] C. McDiarmid and B. A. Reed. Building Heaps Fast. *J. Algorithms*, 10(3):352–365, 1989.
- [16] K. Reinhardt. Sorting *in-place* with a *worst case* complexity of  $n \log n - 1.3n + o(\log n)$  comparisons and  $\epsilon n \log n + o(1)$  transports. In T. Ibaraki et al., editors, *ISAAC*, volume 650 of *LNCS*, pages 489–498. Springer, 1992.
- [17] X.-D. Wang and Y.-J. Wu. An Improved HEAPSORT Algorithm with  $n \lg n - 0.788928n$  Comparisons in the Worst Case. *J. of Comput. Sci. and Techn.*, 22:898–903, 2007.
- [18] I. Wegener. The Worst Case Complexity of McDiarmid and Reed’s Variant of Bottom-Up-Heap Sort is Less Than  $n \lg n + 1.1n$ . In C. Choffrut and M. Jantzen, editors, *STACS*, volume 480 of *LNCS*, pages 137–147. Springer, 1991.
- [19] I. Wegener. BOTTOM-UP-HEAPSORT, a new variant of HEAPSORT, beating, on an average, QUICKSORT (if  $n$  is not very small). *Theor. Comput. Sci.*, 118(1):81–98, 1993.

## APPENDIX

### A Proof of Lem. 3.5

**Proof.** Since the right derivative is monotonically increasing we have:

$$F(x + \delta) - F(x) = \int_x^{x+\delta} F'(t) dt \geq F'(x) \cdot \delta = (\lfloor \lg x \rfloor + 2)\delta$$

and

$$F(y) - F(y - \delta) = \int_{y-\delta}^y F'(t) dt \leq F'(y) \cdot \delta = (\lfloor \lg y \rfloor + 2)\delta$$

This yields:

$$F(y) - F(y - \delta) \leq (\lfloor \lg y \rfloor + 2)\delta \leq (\lfloor \lg x \rfloor + 2)\delta \leq F(x + \delta) - F(x)$$

By adding  $F(x) + F(y - \delta)$  on both sides we obtain the first claim of Lem. 3.5. Note that  $\lim_{\varepsilon \rightarrow 0} F(\varepsilon) = 0$ . Hence the second claim follows from the first by considering the limit  $\delta \rightarrow y$ .  $\square$

### B Proof of Lem. 3.10

**Proof.** First note that the probability for choosing the  $k$ -th element as pivot satisfies

$$\binom{n}{2c+1} \cdot \Pr[\text{pivot} = k] = \binom{k-1}{c} \binom{n-k}{c}.$$

We use the notation of *falling factorial*  $x^\ell = x \cdots (x - \ell + 1)$ . Thus,  $\binom{x}{\ell} = \frac{x^\ell}{\ell!}$ .

$$\begin{aligned} \Pr[\text{pivot} = k] &= \frac{(2c+1)! \cdot (k-1)^c \cdot (n-k)^c}{(c!)^2 \cdot n^{2c+1}} \\ &= \binom{2c}{c} (2c+1) \frac{1}{(n-2c)} \prod_{i=0}^{c-1} \frac{(k-1-i)(n-k-i)}{(n-2i-1)(n-2i)}. \end{aligned}$$

For  $k \leq c$  we have  $\Pr[\text{pivot} = k] = 0$ . So, let  $c < k \leq \frac{n}{2} - \delta n$  and let us consider an index  $i$  in the product with  $0 \leq i < c$ .

$$\begin{aligned} \frac{(k-1-i)(n-k-i)}{(n-2i-1)(n-2i)} &\leq \frac{(k-i)(n-k-i)}{(n-2i)(n-2i)} \\ &= \frac{\left(\left(\frac{n}{2} - i\right) - \left(\frac{n}{2} - k\right)\right) \cdot \left(\left(\frac{n}{2} - i\right) + \left(\frac{n}{2} - k\right)\right)}{(n-2i)^2} \\ &= \frac{\left(\frac{n}{2} - i\right)^2 - \left(\frac{n}{2} - k\right)^2}{(n-2i)^2} \\ &\leq \frac{1}{4} - \frac{\left(\frac{n}{2} - \left(\frac{n}{2} - \delta n\right)\right)^2}{n^2} = \frac{1}{4} - \delta^2. \end{aligned}$$



We have  $\binom{2c}{c} \leq 4^c$ . Since  $2c + 1 \leq \frac{n}{2}$ , we obtain:

$$\Pr[\text{pivot} = k] \leq 4^c(2c + 1) \frac{1}{(n - 2c)} \left( \frac{1}{4} - \delta^2 \right)^c < (2c + 1) \frac{2}{n} \alpha^c.$$

Now, we obtain the desired result.

$$\Pr\left[\text{pivot} \leq \frac{n}{2} - \delta n\right] < \sum_{k=0}^{\lfloor \frac{n}{2} - \delta n \rfloor} (2c + 1) \frac{2}{n} \alpha^c \leq (2c + 1) \alpha^c$$

□

## C More experimental results

We also compare the different pivot selection strategies on the basic QuickHeapsort with no modifications. We test sample of sizes of one, three, approximately  $\lg n$ ,  $\sqrt[4]{n}$ ,  $\sqrt{n/\lg n}$ ,  $\sqrt{n}$ , and  $n^{\frac{3}{4}}$  for the pivot selection.

In Table 2 the average number of comparisons and the standard deviations are listed. We ran the algorithms on arrays of length 10000 and one million. The displayed data is the average resp. standard deviation of 100 runs of QuickHeapsort with the respective pivot selection strategy.

These results are not very surprising: The larger the samples get, the smaller is the standard deviation. The average number of comparisons reaches its minimum with a sample size of approximately  $\sqrt{n}$  elements. One notices that the difference for the average number of comparisons is relatively small, especially between the different pivot selection strategies with non-constant sample sizes. This confirms experimentally that the  $o(n)$ -terms in Thm. 3.1 and Thm. 3.9 are not too big.

$n$	$10^4$		$10^6$	
Sample size	Average number of comparisons	Standard deviation	Average number of comparisons	Standard deviation
1	152573	4.281	21975912	3.452
3	146485	2.169	21327478	1.494
$\sim \lg n$	143669	0.954	20945889	0.525
$\sim \sqrt[4]{n}$	143620	0.857	20880430	0.352
$\sim \sqrt{n/\lg n}$	142634	0.413	20795986	0.315
$\sim \sqrt{n}$	142642	0.305	20783631	0.281
$\sim n^{\frac{3}{4}}$	147134	0.195	20914822	0.168

Table 2: Different strategies for pivot selection tested on  $10^4$  and  $10^6$  elements. The standard deviation of our experiments is given in percent of the average number of comparisons

## D Some Words about the Worst Case Running Time

Obviously the worst case running time depends on how the pivot element is chosen. If just one random element is used as pivot we get the same quadratic worst case running time as for Quicksort. However the probability that in QuickHeapsort we run in such a “bad case” is not higher than in Quicksort, since any choice of pivot elements leading to a worst case scenario in QuickHeapsort also yields the worst case for Quicksort.

If we choose the pivot element as median of approximately  $2 \lg n$  elements, we get a worst case running time of  $\mathcal{O}\left(\frac{n^2}{\lg n}\right)$ , i.e. for the worst case it makes almost no difference, if the pivot is selected as median of  $2 \lg n$  or just as one random element.

However, if we use approximately  $\frac{n}{\lg n}$  elements as sample for the pivot selection, we can get a better bound on the worst case.

Let  $f : \mathbb{N} \rightarrow \mathbb{N}_{\geq 1}$  be some monotonically growing function with  $f \in o(n)$  (e.g.  $f(n) = \lg n$ ). We can apply the ideas of the Median of Medians algorithm [1]: First we choose  $\frac{n}{f(n)}$  random elements, then we group them into groups of five elements each. The median of each group can be determined with six comparisons [13, p. 215]. Now, the median of these medians can be computed using Quickselect. We assume that Quickselect is implemented with the same strategy for pivot selection. That means we get the same recurrence relations for the worst case complexity of the partitioning-phases in QuickHeapsort and for the worst case of Quickselect:

$$T(n) = n + \frac{6n}{5f(n)} + T\left(\frac{n}{5f(n)}\right) + T\left(n - \frac{3n}{10f(n)}\right).$$

This yields  $T(n) \leq cnf(n)$  for some  $c$  large enough. Hence with this pivot selection strategy, we reach a worst case running time for QuickHeapsort of  $n \lg n + \mathcal{O}(nf(n))$  and – if  $f(n) \in \omega(1)$  – average running time as stated in Sect. 3.

Driving this strategy to the end and choosing  $f(n) = 1$  leads to Ultimate Heapsort (or better a slight modification of it – and Quickselect turns into the Median of Medians algorithm). Then we have  $T(n) = n \lg n + \mathcal{O}(n)$  for the worst case of QuickHeapsort. However, our bound for the average case does not hold anymore.

In order to obtain an  $n \lg n + \mathcal{O}(n)$ -bound for the worst case without losing our bound for the average case, we can apply a simple trick: Whenever after the partitioning it turns out that the pivot does not lie in the interval  $\{\frac{n}{4}, \dots, \frac{3n}{4}\}$  we switch to Ultimate Heapsort. This immediately yields the worst case bound of  $n \lg n + \mathcal{O}(n)$ . Moreover, the proof of Thm. 3.9 can easily be changed in order to deal with this modification: Let  $C \cdot n$  be the worst case number of comparisons for pivot selection and partitioning in Ultimate Heapsort. We can change Eq. (8) to

$$\Pr[\text{pivot} < \mu - \delta n] \leq \frac{1}{8C}\epsilon.$$

Then, the rest of the proof is exactly the same. Hence, Thm. 3.9 and Thm. 3.1 are also valid when switching to Ultimate Heapsort in the case of a ‘bad’ choice of the pivot.

## E Pseudocode of basic QuickHeapsort

---

**Algorithm 5.1**

---

```
procedure QuickHeapsort( $A[1..n]$ )
begin
  if  $n > 1$  then
     $p := \text{ChoosePivot};$ 
     $k := \text{PartitionReverse}(A[1..n], p);$ 
    if  $k \leq n/2$  then
      TwoLayerMaxHeap( $A[1..n], k - 1$ );           (* heap-area:  $\{1..k-1\}$  *)
      swap( $A[k], A[n - k + 1]$ );
      QuickHeapsort( $A[1..n - k]$ );                 (* recursion *)
    else
      TwoLayerMinHeap( $A[1..n], n - k$ );           (* heap-area:  $\{k+1..n\}$  *)
      swap( $A[k], A[n - k + 1]$ );
      QuickHeapsort( $A[(n - k + 2)..n]$ );           (* recursion *)
    endif
  endif
endprocedure
```

---

The ChoosePivot function returns an element  $p$  of the array chosen as pivot. The PartitionReverse function returns an index  $k$  and rearranges the array  $A$  so that  $p = A[k]$ ,  $A[i] \geq A[k]$  for  $i < k$  and  $A[i] \leq A[k]$  for  $i > k$  using  $n - 1$  comparisons.

---

**Algorithm 5.2**

---

```
function SpecialLeaf( $A[1..m]$ ):
begin
   $i := 1$ ;
  while  $2i \leq m$  do                               (* i.e. while  $i$  is not a leaf *)
    if  $2i + 1 \leq m$  and  $A[2i + 1] > A[2i]$  then
       $A[i] := A[2i + 1];$ 
       $i := 2i + 1$ ;
    else
       $A[i] := A[2i];$ 
       $i := 2i$ ;
    endif
  endwhile
  return  $i$ ;
endfunction
```

---

---

**Algorithm 5.3**

---

```
procedure TwoLayerMaxHeap( $A[1..n], m$ )
begin
  ConstructHeap( $A[1..m]$ );
  for  $i := 1$  to  $m$  do
```

```
    temp :=  $A[n - i + 1]$ ;  
     $j := \text{SpecialLeaf}(A[1..m])$ ;  
     $A[j] := temp$ ;  
  endfor  
endprocedure
```

---

The procedure `TwoLayerMinHeap` is symmetric to `TwoLayerMaxHeap`, so we do not present its pseudocode here.